

Forecasting Stock Price Using Machine Learning Technique

Mohammed Nazim Uddin, PhD
School of Science, Engineering and Technology
East Delta University
Chattagram, Bangladesh
nazim@eastdelta.edu.bd

Md Tanvir Rahman
School of Science, Engineering and Technology
East Delta University
Chattagram, Bangladesh
Ornobtanvir.git@gmail.com

Abstract—Stock market is an emerging sector in any country of the world. Many people directly related to this sector. Stock market prediction is the act of trying to determine the future value of a company stock or other financial instrument. When publicly traded companies issue shares of stock to investors, each of those shares is assigned a monetary value, or price. Stock prices can go up or down depending on different factors. Stock prices can be affected by a number of things including volatility in the market, current economic conditions, and popularity of the company. The successful prediction of a stock's future price could yield significant profit. Along with the development with the stock market, forecasting become an important topic. Since finance market has become more and more competitive, stock price prediction has been a hot research topic in the past few decades. Predicting stock price is regarded a challenging task because stock market is essentially non linear, non-parametric, noisy, and a chaotic system. Trend of a market depends on many things like liquid money human behavior, news related to stock market etc. All this together controls the behavior of trends in a stock market with the advancement of the computing technology we use machine learning technique, like Support Vector Regression, K-nearest-neighbor, linear Regression, Random forest Regressor, for analyzing time series data to predict stock price. In this paper we try to develop a forecasting model with stacking multiple method to find the best forecast of the stock price.

Keywords— Time Series data, SVR, KNN-Regressor, liner regression, stacking regressor, Random forest regression

I. INTRODUCTION

The goal is to take time series data, find the equation that best fits the data, and be able to forecast out a specific value. Time series data is a continuous data statistical observations recorded over a specific period of time. This model will try to understand the pattern of the continuous data[1] by combining different methods and produce a best fit line that fits the data. The target is to determine the future stock price and improve their strategy for future. regression[2] models are among the most known regression models used in the machine learning community and recently many researchers have examined their sufficiency in ensembles[3][4]. Although many methods of ensemble design have been proposed, there is as yet no obvious picture of which method is best. One notable successful adoption of ensemble learning is the distributed scenario. In this work, we propose an efficient distributed method that uses the same training set with the parallel usage of an averaging methodology that combines linear regression[5] and KNN regression

models[6], Support Vector Regression[7], random Forest Regression[8][9]. We performed a comparison of the presented ensemble with other ensembles that use either the linear regression as base learner and the performance of the proposed method was better in most cases. Using averaging methodology, we expect to obtain better results because both theory and experiments show that averaging helps most if the errors in the individual regression models are not positively correlated. linear regression is a linear approach to modelling the relationship between a scalar response dependent variable and one or more explanatory variables (or independent variables). The case of one explanatory variable is called simple linear regression. For more than one explanatory variable, the process is called multiple linear regression. support vector machines are supervised learning models with associated learning algorithms that analyze data used for classification and regression analysis. Given a set of training examples, each marked as belonging to one or the other of two categories, an SVM training algorithm builds a model by creating a hyper plane[9] that assigns new examples to one category or the other. In support vector regression (SVR). The model produced by support vector classification (as described above) depends only on a subset of the training[10] data, because the cost function for building the model does not care about training points that lie beyond the margin. In pattern recognition, the k-nearest neighbors algorithm (k-NN) is a non-parametric method used for classification and regression. In both cases, the input consists of the k closest training examples in the feature space. The output depends on whether KNN is used for classification or regression. In case of the kNN regression the output is the property value for the object. This value is the average of the values of its k nearest neighbors. a useful technique can be used to assign weight to the contributions of the neighbors, so that the nearer neighbors contribute more to the average than the more distant ones. A model that combines KNN regression, Linear regression, Support Vector Regression, Random Forest regression model used for predicting stock prices can forecast better price accuracy. The paper is structured as follows. Section II presents the most well-known methods for building ensembles. Section III contains Proposed methodology. Section IV contains Architecture Overview, Section V contains Implementation and Evaluation. Section VI contains conclusion.

II. ENSEMBLES OF REGRESSION MODELS

Bagging is a "bootstrap" ensemble method that creates individualized by randomly drawing, with replacement, N examples - where N is the size of the original set; many of

the original examples may be repeated in the resulting training set while others may be left out. After construction of several regression models, averaging the

predictions of each regression model performs the final prediction. Instability (responsiveness to changes in the training data) is a prerequisite for bagging to be effective. Another approach for building ensembles of regression models is to use a variety of learning algorithms on all of the training data and combine their predictions. When multiple regression models are combined using averaging methodology, we expect to obtain good results based on the belief that the majority of experts are more likely to be correct in their decision when they are close in their opinions. Stacked generalization or Stacking, is a more sophisticated approach for combining predictions of different learning algorithms. Stacking combines multiple regression models to induce a higher-level regression model with improved performance. In detail, the original data set constitutes the level zero data and all the base regression models run at this level. The level one data are the outputs of the base regression models. A learning algorithm is then used to determine how the individuals for its ensemble by training each regression model on a random redistribution of the training set. Each regression model's training set is generated by randomly drawing, with replacement, N examples - where N is the size of the original set; many of the accurately learned. After several cycles, the prediction is performed by taking a weighted average of the predictions of each regression model, with the weights being proportional to each regression model's performance on its training set.

III. PROPOSED METHODOLOGY

Generate linear regression algorithm from this formula

$$m = (N \sum xy - \sum x \sum y) / (N \sum x^2 - (\sum x)^2) \quad (1)$$

$$b = \sum y - m(\sum x) / N \quad (2)$$

$$y = mx + b \quad (3)$$

Generate knn-regression by computing euclidean distance from the query example to the labeled example. and ordering the example by increasing rate .we estimate euclidean distance by following formula

$$\text{dist}(A, B) = \sqrt{\sum (x_i - y_i)^2 / 2m} \quad (4)$$

generate SVM regression, the input (X) is first mapped onto a m -dimensional feature space using some fixed (nonlinear) mapping, and then a linear model is constructed in this feature space. Using $f(x, w)$ mathematical notation, the linear model (in the feature space) is given by

$$f(\mathbf{x}, \mathbf{w}) = \sum_{j=1}^m w_j g_j(\mathbf{x}) + b$$

loss function is estimated by this formula

$$L_\epsilon(y, f(\mathbf{x}, \mathbf{w})) = \begin{cases} 0 & \text{if } |y - f(\mathbf{x}, \mathbf{w})| \leq \epsilon \\ |y - f(\mathbf{x}, \mathbf{w})| - \epsilon & \text{otherwise} \end{cases}$$

Then stacking is performed using these algorithm. Stacking is concerned with combining multiple classifiers generated by different learning algorithms L_1, \dots, L_n on a single data set S , which is composed by a feature vector $S^i = (X^i, Y^i)$.

- The stacking process can be broken into two phases:

1. Generate a set of base-level classifiers C_1, \dots, C_n .

- Where $C_i = L_i(S)$

2. Train a meta-level classifier to combine base level classifier

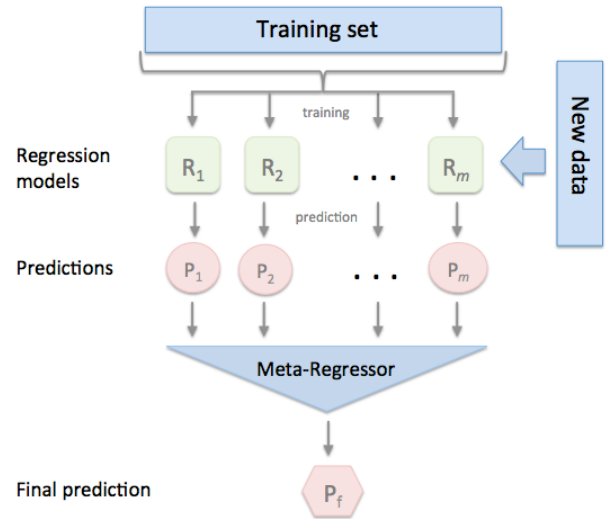
3. The training set for the meta-level classifier is generated through a leave-one-out cross validation process.

$$\forall i = 1, \dots, n \text{ and } \forall k = 1, \dots, N$$

$$C_k = L_k(S - S^i)$$

- The learned classifiers are then used to generate predictions for $Y_k^i = C_k(x^i)$

- The meta-level data sets consists of examples of the form where the features are the predictions of the base-level classifiers and the class is the correct class of the example in hand.



IV. ARCHITECTURE OVERVIEW

The prediction system has a two tier architecture top tier is dedicated to preparing the data sets from multiple information sources to make them ready for the predication tasks in the next tier. It is composed of two major parts. The first part is data preprocessing. In this process we process the data by adding more feature and removing unnecessary feature and removing the bad data and also the absence of the data. The second part is the data alignment. The second tier is dedicated to the market volatility analysis and prediction through the model integration and training, which uses multiple kernel learning methodology to train the model. It consists of three tasks: First, we build one regression model per source. Second, we train the model with the same data sets, then we create a stacked algorithm using these algorithm. In this paper, we use the multi-kernel learning method.

A. Data Preprocessing

In the real world, many data sets are very messy. Most stock price/volume data is pretty clean, rarely with missing data, but many data sets will have a lot of missing data. filter out other unimportant feature from the feature because not all the feature will be included into the final feature list. The reason behind it is the unnecessary feature and those value which has no relation with the stock market prediction will reduce the accuracy of the prediction.

B. Feature Extraction

We only work with the valuable feature, created some new valuable feature through manipulation. Right feature can help to get more accurate results.

C. Training The ML Model

The way this works is taking 75% of data, and use this to train the machine learning classifier. Then taking the remaining 25% of our data, and test the classifier. Since this is sample data, we should have the features and known labels. Thus, if we test on the last 25% of our data, we can get a sort of accuracy and reliability

D. Creating Stacked algorithm

Create the stacked algorithm using the model previously made. And averaging the prediction for creating the final prediction.

V. IMPLEMENTATION AND EVALUATION

A. Data Collection

Data collection is the process of gathering and measuring information on targeted variables in an established systematic fashion. Which enables us to evaluate outcome. In this case for regression we used stock price data set of Google, Apple, Microsoft, IBM, Nike which is collected from Quandl which is a platform for financial data.

VI. RESULT ANALYSIS

TABLE I. ACCURACIES OF THE PROPOSED MODEL

Company	Accuracy Based on Cross Validation				
	Line ar regr essi onm	Knn regres sion	SVR	RFG	Proposed model
Google	87.97%	89.18%	74.32%	88.67%	90.72%
Apple	92.2%	91.4%	91.0%	93.4%	94.8%
Microsoft	91.4%	91.4%	91.4%	91.4%	91.4%
IBM	73.6%	71.3%	77.8%	78.4%	80.8%
Nike	94.08%	94.9%	95.05%	92.9%	96.8%

After determining and comparing with other models. In our proposed model we have attained the highest accuracy among all others.

VII. CONCLUSIONS AND FUTURE WORK

It is known that if we are only concerned for the best possible correlation coefficient, it might be difficult or impossible to find a single regression model that performs as well as a good ensemble of regression models. In this study, we built an ensemble of regression models using four different learning methods.

REFERENCES

- [1] Classification with discrete and continuous variables via general mixed-data models, Alexander de Leon: University of Calgary, Andrea Soo: University of Calgary
- [2] J. Fox, Applied Regression Analysis, Linear Models, and Related Methods, ISBN: 080394540X, Sage Pubns (1997).
- [3] L. Breiman, Bagging Predictors. Machine Learning, 24(3) (1996) 123-140.
- [4] D. Opitz, R. Maclin, Popular Ensemble Methods: An Empirical Study, Artificial Intelligence Research, 11 (1999): 169-198, Morgan Kaufmann. 12.C. Perhch, F. J. Provost, J. S.
- [5] Linear regression Analysis on Net Income of an Agronomical Company, Supichaya Sunthornjittanon: Portland State University
- [6] Efficient and Accurate knn based Classification and Regression, Harshit Dubey: International Institute of Information Technology
- [7] Design and Training of Support Vector Machines, Alistair Shilton: The University of Melbourne
- [8] Advances in Random Forests with Application to Classification, Arnu Pretorius: Stellenbosch University
- [9] Cristianini and J. Shawe-Taylor. An Introduction to Support Vector Machines and other kernel-based learning method. Cambridge University Press, Cambridge, UK, 2000
- [10] L. Breiman, Stacked Regression. Machine Learning, 24 (1996): 49-64. T.G. Dietterich, Ensemble methods in machine learning. In Kittler, J., Roli, F., eds.:
- [11] Multiple Classifier Systems. LNCS Vol. 1857, Springer (2001) 1-1 N. Duffy, D. Helmbold, Boosting Methods for Regression, Machine Learning, 47, (2002)
- [12] J. Friedman, Stochastic Gradient Boosting, Computational Statistics and Data Analysis
- [13] Y. Grandvalet, Bagging Equalizes Influence, Machine Learning, Volume 55(3) (2004)
- [14] N.L. Hjort, G. Claeskens, Frequentist Model Average Estimators, Journal of the American Statistical Association,
- [15] Y. Morimoto, H. Ishii, S. Morishita, Efficient Construction of Regression Trees with Range and Region Splitting, Machine Learning
- [16]
- [17] Simonoff, Tree Induction vs. Logistic Regression: A Learning-Curve Analysis. Journal of Machine Learning Research 4 (2003) 211-255
- [18] Y. Wang, I. H. Witten, Induction of model trees for predicting continuous classes, In Proc. of the Poster Papers of the European Conference on ML, (1997) 128-137.